

Information capacity in the weak-signal approximation

Lubomir Kostal*

Institute of Physiology AS CR, v.v.i., Videnska 1083, 142 20 Prague 4, Czech Republic

(Dated: August 10, 2010: Accepted for publication in Physical Review E)

We derive an approximate expression for mutual information in a broad class of discrete-time stationary channels with continuous input, under the constraint of vanishing input amplitude or power. The approximation describes the input by its covariance matrix, while the channel properties are described by the Fisher information matrix. This separation of input and channel properties allows us to analyze the optimality conditions in a convenient way. We show that input correlations in memoryless channels do not affect channel capacity since their effect decreases fast with vanishing input amplitude or power. On the other hand, for channels with memory, properly matching the input covariances to the dependence structure of the noise may lead to almost noiseless information transfer, even for intermediate values of the noise correlations. Since many model systems described in mathematical neuroscience and biophysics operate in the high noise regime and weak-signal conditions, we believe, that the described results are of potential interest also to researchers in these areas.

PACS numbers: 89.90.+n, 89.70.Kn

I. INTRODUCTION

Information theory is a mathematical framework that provides tools for quantification of information content and information transfer in systems defined by general probabilistic rules [1]. The theory has been applied successfully to a wide range of problems [2], including, e.g., classical and quantum computation and communication [3–5], optical communication [6–8] or quantification of different aspects of information processing in real neurons and neuronal models [9–15].

The measure of information transfer in information theory is represented by a nonlinear functional of the probability measure over the joint input-output space [1]. The concavity of this functional in the input probability measure has important implications for numerical approaches to finding the information optimality conditions [1, 16–18]. On the other hand, approximations or even closed-form solutions are quite rare. The classical exact solution for the linear channel with additive (possibly non-white) Gaussian noise [1, 19] and input power constraint has been applied in many different situations. However, in many cases of interest the channel is significantly nonlinear or non-Gaussian or there are different input constraints [20] and one has to rely on numerical solutions or approximations.

The approximations allow us to investigate, although locally and under perhaps restrictive scenario, the effect of individual components in the system on the optimality conditions. In particular, if the noise in information transfer is substantially low and regular, there exists a tight lower bound on the information optimality conditions (denoted as *low-noise* approximation in this paper) which has been investigated in [12, 21–23]. In this paper we continue the effort started in [24] and we describe

essentially the opposite situation: the *high-noise* approximation. Such approximation is of interest when the signal is very weak compared to the noise in the information transfer, for example, as in the classical stochastic resonance effect observed in electrosensory neurons [24, 25].

II. MEASURES OF INFORMATION

Throughout this paper we assume the discrete-time setting [5], we denote the consequent channel outputs (responses) as a vector of random variables (r.v.) $R = (\{R_i\}_{i=1}^n)^T$, which may be discrete or continuous, i indexes the time and $(\cdot)^T$ denotes the transposition. The response, $R_i = r_i$, results from the corresponding input $\Theta_i = \theta_i$, where the input is also described by a n -dimensional r.v. Θ . The multidimensional description of the process of information transfer between Θ and R allows us to include the effect of memory, i.e., the dependence on current and also on past inputs and responses. We also assume that the input alphabet is continuous [5]. In the following we consider stationary channels fully described by the conditional probability density function (p.d.f.) $f(\mathbf{r}|\boldsymbol{\theta})$, which generally factorizes as [26]

$$f(\mathbf{r}|\boldsymbol{\theta}) = \prod_{i=1}^n f_i(r_i|\theta_i, \theta_{i-1}, \dots, \theta_1, r_{i-1}, \dots, r_1). \quad (1)$$

We do not consider channel feedback, the dependence of current input on past responses [1].

Mutual information (MI) is the fundamental quantity measuring information transfer in channels [1]. MI $I(\Theta; R)$ gives the degree of statistical dependence between inputs and responses, defined as

$$I(\Theta; R) = \langle D_{\text{KL}} [f(\mathbf{r}|\boldsymbol{\theta}) \| p(\mathbf{r})] \rangle_{\boldsymbol{\theta}}, \quad (2)$$

where

$$p(\mathbf{r}) = \langle f(\mathbf{r}|\boldsymbol{\theta}) \rangle_{\boldsymbol{\theta}} \quad (3)$$

* kostal@biomed.cas.cz

is the marginal joint p.d.f. of responses, and the averaging is with respect to the input p.d.f., $\pi(\boldsymbol{\theta})$. The Kullback-Leibler (KL) divergence is defined as

$$D_{\text{KL}}[f(\mathbf{r}|\boldsymbol{\theta}) \parallel p(\mathbf{r})] = \left\langle \ln \frac{f(\mathbf{r}|\boldsymbol{\theta})}{p(\mathbf{r})} \right\rangle_{\mathbf{r}|\boldsymbol{\theta}}, \quad (4)$$

where the averaging is with respect to $f(\mathbf{r}|\boldsymbol{\theta})$. From Eq. (2) follows, that MI is a property of the joint distribution of stimuli and responses. Of particular interest are the *optimality conditions* for information transfer, that is the maximum value of $I(\boldsymbol{\Theta}; \mathbf{R})$ and the corresponding optimal $\pi(\boldsymbol{\theta})$. In order to have a well-posed problem, one is interested in the optimality conditions for $\boldsymbol{\Theta}$ satisfying certain additional constraints, e.g., average power or range of inputs [1, 20]. The maximum value of MI per channel use, taken over all possible stimuli distributions satisfying constraints \mathcal{G} , is denoted as the information capacity, \mathcal{C} , defined as [20]

$$\mathcal{C} = \lim_{n \rightarrow \infty} \frac{1}{n} \left[\sup_{\pi(\boldsymbol{\theta}) \in \mathcal{G}} I(\boldsymbol{\Theta}; \mathbf{R}) \right]. \quad (5)$$

In this paper we interpret \mathcal{C} as the upper bound on the rate at which the information can be transmitted reliably [1], without considering the complexity of achieving such maximum rate in practical terms. Specifically, do not discuss the properties of any particular coding and decoding schemes [5].

Whenever we are interested in reliability of input-output transmission, we naturally interfere with the domain of statistical estimation theory [27]. Fisher information (FI) matrix, defined as

$$\mathbf{J}(\boldsymbol{\theta}|\mathbf{R}) = \langle [\nabla \ln f(\mathbf{r}|\boldsymbol{\theta})][\nabla \ln f(\mathbf{r}|\boldsymbol{\theta})]^\top \rangle_{\mathbf{r}|\boldsymbol{\theta}}, \quad (6)$$

where

$$\nabla = \left(\frac{\partial}{\partial \theta_1}, \dots, \frac{\partial}{\partial \theta_n} \right)^\top, \quad (7)$$

imposes limits on the precision of $\boldsymbol{\theta}$ estimation from the responses by means of the Cramer-Rao bound, which says that for the variance of any unbiased estimator of θ_i holds $\text{Var}(\hat{\theta}_i) \geq [\mathbf{J}^{-1}(\boldsymbol{\theta}|\mathbf{R})]_{ii}$ [27]. Generally, FI requires that $f(\mathbf{r}|\boldsymbol{\theta})$ is continuously differentiable in $\boldsymbol{\theta}$ [27]. In this paper, we additionally assume that $f(\mathbf{r}|\boldsymbol{\theta})$ is twice continuously differentiable in $\boldsymbol{\theta}$, so that the following conditions hold

$$\int_{\mathbf{R}} \nabla f(\mathbf{r}|\boldsymbol{\theta}) d\mathbf{r} = \mathbf{0}, \quad \int_{\mathbf{R}} \nabla \nabla^\top f(\mathbf{r}|\boldsymbol{\theta}) d\mathbf{r} = \mathbf{0}. \quad (8)$$

There is a variety of relationships between FI, MI and KL divergence established in the literature [1, 28, 29], further motivated by the fields of information geometry [30] or stochastic complexity [31]. The already mentioned *low-noise* approximation to MI is constructed by employing the Cramer-Rao bound [12, 21–23]. Although we demonstrate that the *high-noise* approximation also involves FI, we never employ the Cramer-Rao bound and the appearance of FI is due to certain asymptotic properties of the KL distance [28].

III. INFORMATION TRANSFER BY WEAK SIGNALS

A. Small input amplitude limit

The channel properties are described by the conditional probability density $f(\mathbf{r}|\boldsymbol{\theta})$, which satisfies the regularity conditions (8). The input, described by r.v. $\boldsymbol{\Theta}$, is restricted in amplitude,

$$\boldsymbol{\Theta} \in [\boldsymbol{\theta}_0 - \Delta\boldsymbol{\theta}, \boldsymbol{\theta}_0 + \Delta\boldsymbol{\theta}], \quad (9)$$

for chosen $\boldsymbol{\theta}_0$ and $\Delta\boldsymbol{\theta}$, or more precisely in components: for all i holds $\Theta_i \in [\theta_{0i} - \Delta\theta_i, \theta_{0i} + \Delta\theta_i]$ and $\Delta\theta_i > 0$. The situation for a memoryless channel is illustrated in Fig. 1. The goal is to derive an approximation to mutual information in the limit $\|\Delta\boldsymbol{\theta}\| \rightarrow 0$. We demonstrate in detail in Appendix A, that the approximation (to second order in the input amplitude) can be written as

$$I(\boldsymbol{\Theta}; \mathbf{R}) \approx \frac{1}{2} \text{tr} [\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \mathbf{C}_{\boldsymbol{\Theta}}], \quad (10)$$

where $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})$ is the FI matrix from Eq. (6) evaluated at $\boldsymbol{\theta} = \boldsymbol{\theta}_0$, $\mathbf{C}_{\boldsymbol{\Theta}}$ is the covariance matrix of $\boldsymbol{\Theta}$ and $\text{tr}(\cdot)$ is the matrix trace. Eq. (10), derived also in [24], holds for a broad class of channels with memory, both biologically-inspired and artificial and represents the main result. An important feature of Eq. (10) is, that the channel properties (described by the FI matrix) and the input properties (described by its covariance matrix) are separated. Therefore, the maximum value of MI can be found by matching the corresponding elements of $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})$ and $\mathbf{C}_{\boldsymbol{\Theta}}$. The elements of the covariance matrix of $\boldsymbol{\Theta}$ can be written as [32]

$$[\mathbf{C}_{\boldsymbol{\Theta}}]_{ik} = \sigma^2 \varrho_{ik}, \quad (11)$$

where $\sigma^2 \equiv \sqrt{\text{Var}(\Theta_i)\text{Var}(\Theta_k)}$ is constant for all i, k due to stationarity, and $\varrho_{ik} = \text{corr}(\Theta_i, \Theta_k)$ is the correlation coefficient. The maximum variance of the amplitude constrained input from Eq. (9) is $\max \sigma^2 = (\Delta\theta)^2$ and $-1 < \varrho_{ik} < 1$, thus $I(\boldsymbol{\Theta}; \mathbf{R})$ in Eq. (10) is maximized if

$$\varrho_{ik} \rightarrow \text{sgn}[\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})]_{ik}, \quad (12)$$

where $\text{sgn}(\cdot)$ is the sign function. Note, that the diagonal elements of the FI matrix are positive while the off-diagonal elements can be negative. It may happen, that the matrix $\mathbf{C}_{\boldsymbol{\Theta}}$ formed by Eqns. (12) and (11) is not positive-semidefinite [33], i.e., it cannot be a proper covariance matrix [34], even though $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})$ generally is positive-semidefinite [27]. However, in all problems we have calculated so far, proper input covariance matrix could be formed, given $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})$, and then it holds from Eqns. (5) and (10)

$$\mathcal{C} \approx \mathcal{C}_{\text{high}} = \lim_{n \rightarrow \infty} \frac{(\Delta\theta)^2}{2n} \sum_{i,k} |[\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})]_{ik}|, \quad (13)$$

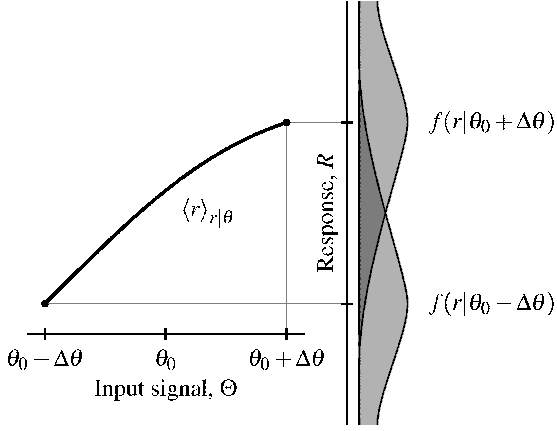


FIG. 1. Information transmission with amplitude-constrained inputs. The input signal, described by r.v. Θ , is restricted to the interval $[\theta_0 - \Delta\theta, \theta_0 + \Delta\theta]$. Due to presence of noise, the responses to each particular θ vary randomly, described by the conditional probability density $f(r|\theta)$. While the memoryless information channel is fully described by $f(r|\theta)$, the amount of information transferred depends on both $f(r|\theta)$ and the distribution of Θ . We examine the maximum information transfer by inputs restricted to small amplitudes when there is a significant overlap of $f(r|\theta_0 - \Delta\theta)$ and $f(r|\theta_0 + \Delta\theta)$. Heuristically, the problem can be also described as the information transmission in a very noisy environment, or under very low signal-to-noise ratio conditions.

where $\mathcal{C}_{\text{high}}$ denotes the *high noise* approximation to the true capacity \mathcal{C} .

For stationary memoryless channels $f(\mathbf{r}|\boldsymbol{\theta})$ factorizes due to Eq. (1) as [1, p.193]

$$f(\mathbf{r}|\boldsymbol{\theta}) = \prod_{i=1}^n f(r_i|\theta_i), \quad (14)$$

thus from Eq. (6) follows that the FI matrix is diagonal, $J(\theta_0|R) \equiv [\mathbf{J}(\theta_0|\mathbf{R})]_{ii} = \langle [\partial_\theta \ln f(r|\theta)]^2 \rangle_{r|\theta}$, and from Eq. (13) we have

$$\mathcal{C}_{\text{high}} = \frac{(\Delta\theta)^2}{2} J(\theta_0|R), \quad (15)$$

a result obtained by different means in [35]. The optimal input p.d.f., $\pi^*(\theta)$, is the maximum variance distribution over the given input range,

$$\pi^*(\theta) = \frac{1}{2} \delta(\theta - \theta_0 - \Delta\theta) + \frac{1}{2} \delta(\theta - \theta_0 + \Delta\theta), \quad (16)$$

where $\delta(\cdot)$ is the Dirac's delta function. In other words, the capacity is achieved by a binary input, and thus $\mathcal{C} \leq 1$ bit.

From Eq. (10) follows, that non-diagonal elements of \mathbf{C}_Θ do not affect the information capacity of memoryless channels in the vanishing input amplitude case. This result is counterintuitive, because correlations generally

decrease the input entropy [1]. Therefore in the following we provide a proof which is independent of Eq. (10). Let us consider two consequent uses of a stationary memoryless channel, i.e., $\boldsymbol{\Theta} = \{\Theta_1, \Theta_2\}^\top$, $\mathbf{R} = \{R_1, R_2\}^\top$. We assume, that the inputs Θ_1 and Θ_2 are generally statistically dependent, $(\Theta_1, \Theta_2) \sim \pi(\theta_1, \theta_2)$, and the joint marginal distribution of responses is denoted as $p(\mathbf{r})$, see also Eq. (3). By employing the factorization (14) and basic relations between entropy, $h(\mathbf{R}) = -\langle \ln p(\mathbf{r}) \rangle_{\mathbf{r}}$, and MI [1, p.21] we have

$$\begin{aligned} I(\boldsymbol{\Theta}; \mathbf{R}) &= h(\mathbf{R}) - \langle h(\mathbf{R}|\boldsymbol{\theta}) \rangle_{\boldsymbol{\theta}} = \\ &= h(R_1) + h(R_2) - I(R_1; R_2) - \\ &\quad - \langle h(R_1|\theta_1) + h(R_2|\theta_2) \rangle_{\boldsymbol{\theta}} = \\ &= I(\Theta_1; R_1) + I(\Theta_2; R_2) - I(R_1; R_2) = \\ &= 2I(\Theta_1; R_1) - I(R_1; R_2), \end{aligned} \quad (17)$$

since $I(\Theta_1; R_1) = I(\Theta_2; R_2)$ due to stationarity. In other words, the difference in information transfer when using two dependent or independent inputs in the memoryless channel is equal to $I(R_1; R_2)$. Obviously, for Θ_1, Θ_2 independent holds $I(R_1; R_2) = 0$. The strength of the dependence between R_1 and R_2 for correlated inputs depends on the input range and the conditional response distributions, see Fig. 1. We expect $I(R_1; R_2)$ to be maximal for the extreme input dependence, e.g., $\Theta_2 = \Theta_1$, where Θ_1 is equiprobably equal either to $\theta_0 - \Delta\theta$ or $\theta_0 + \Delta\theta$. It follows, that R_1, R_2 are conditionally (given Θ_1) identically and conditionally independently distributed. If $f(r|\theta_0 - \Delta\theta)$ and $f(r|\theta_0 + \Delta\theta)$ are well separated, then $I(R_1; R_2) > 0$ because R_2 provides redundant information to R_1 . As $\Delta\theta \rightarrow 0$, then $f(r|\theta_0 - \Delta\theta)$ and $f(r|\theta_0 + \Delta\theta)$ become (almost) identical due to continuity in θ and thus $I(R_1; R_2) \rightarrow 0$. To make the argument precise, we show that $I(R_1; R_2) = 0$ to the second order in the input amplitude, so that the effect of input correlations in memoryless channels is of higher order than the approximate Eq. (10). The joint response distribution is

$$\begin{aligned} p(r_1, r_2) &= \frac{1}{2} f(r_1|\theta_0 + \Delta\theta) f(r_2|\theta_0 + \Delta\theta) + \\ &\quad + \frac{1}{2} f(r_1|\theta_0 - \Delta\theta) f(r_2|\theta_0 - \Delta\theta), \end{aligned} \quad (18)$$

from which the marginals follow $p(r_1) = f(r_1|\theta_0 + \Delta\theta)/2 + f(r_1|\theta_0 - \Delta\theta)/2$, and similarly for $p(r_2)$. We employ another formula for MI [1, p.251]

$$I(R_1; R_2) = D_{\text{KL}}[p(r_1, r_2) \parallel p(r_1)p(r_2)]. \quad (19)$$

By substituting from Eq. (18) into Eq. (19), and by employing the Taylor expansion in $\Delta\theta$ around $\Delta\theta = 0$, we have (the terms up to $\Delta\theta$ are zero)

$$I(R_1; R_2) \approx (\Delta\theta)^2 \iint_{R_1 \times R_2} \left[\frac{\partial f(r_1|\theta)}{\partial \theta} \frac{\partial f(r_2|\theta)}{\partial \theta} \right] \Big|_{\theta=\theta_0} dr_1 dr_2,$$

which is equal to zero, due to Eq. (8). The first nonzero term is of 4-th order, and can be written as

$(\Delta\theta)^4 J(\theta_0|R_1)J(\theta_0|R_2)/2$, provided that $f(r|\theta)$ is three times continuously differentiable in θ .

On the other hand, for channels with memory the input correlations do matter, irrespectively of the smallness of the amplitude. Consider, for example, two channel uses in the additive noise case, $R_i = \Theta_i + Z_i$, $\langle Z_i \rangle = 0$, where $i = 1, 2$. It is possible to approach the noiseless channel in the extreme case of matching input and noise correlations in accord with Eq. (10), e.g., if $\text{corr}(Z_1, Z_2) \rightarrow -1$ and $\text{corr}(\Theta_1, \Theta_2) \rightarrow 1$, then $R_1 = \Theta_1 + Z_1$ and $R_2 = \Theta_1 - Z_1$ and so by adding $R_1 + R_2$ we can recover the value of Θ_1 perfectly.

B. Small input power limit

The signal power [36], P_Θ , of an input signal described by r.v. Θ is defined as

$$P_\Theta = \frac{1}{n} \langle \Theta^\top \Theta \rangle. \quad (20)$$

For the covariance matrix \mathbf{C}_Θ of r.v. Θ holds $\mathbf{C}_\Theta = \langle (\Theta - \langle \Theta \rangle)(\Theta - \langle \Theta \rangle)^\top \rangle$, and therefore

$$P_\Theta = \frac{1}{n} [\text{tr} \mathbf{C}_\Theta + \|\langle \Theta \rangle\|^2]. \quad (21)$$

The information channel is constrained in the input power P if only inputs that satisfy $P \geq P_\Theta$ are considered. It is common in information theory of power-constrained channels, to assume $\langle \Theta \rangle = \mathbf{0}$, then $P_\Theta = \text{tr} \mathbf{C}_\Theta / n$ [1, p.277], which we assume here also. The assumption $\langle \Theta \rangle = \mathbf{0}$ results in simpler notation, although it does not affect the generality of results. Due to stationarity, the marginal variances of r.v. Θ are constant, $\text{Var}(\Theta_i) = \text{const.}$ for all i , thus we can write

$$\Theta = \varepsilon \tilde{\Theta}, \quad (22)$$

where $\text{Var}(\tilde{\Theta}_i) = 1$ and $\varepsilon > 0$ is the scaling factor. The power of the input is then $P_\Theta = \varepsilon^2$, and the vanishing input power is achieved by $\varepsilon \rightarrow 0$.

The approximate expression for MI in the vanishing input power limit is obtained analogously to the proof presented in Appendix A, by expressing $I(\Theta; \mathbf{R})$ in terms of the auxiliary r.v. $\tilde{\Theta}$, and then expanding for $\varepsilon \rightarrow 0$ around $\varepsilon = 0$. Let $\Theta \sim \pi(\theta)$ and $\tilde{\Theta} \sim g(\tilde{\theta})$, then from Eq. (22) follows $\pi(\theta) = g(\tilde{\theta}/\varepsilon)/\varepsilon = g(\tilde{\theta})/\varepsilon$, and also $d\theta = \varepsilon d\tilde{\theta}$. The MI can be written by (analogously to Eq. (A2))

$$I(\Theta; \mathbf{R}) = \langle D_{\text{KL}} [f(\mathbf{r}|\varepsilon\tilde{\theta}) \| \langle f(\mathbf{r}|\varepsilon\tilde{\theta}) \rangle_{\tilde{\theta}}] \rangle_{\tilde{\theta}}. \quad (23)$$

The rest follows the argument of Appendix A, although simplified due to $\langle \Theta \rangle = \mathbf{0}$. It is obvious from the general proof, that the assumption on zero $\langle \Theta \rangle$ is not essential, only that the vanishing input power is then with respect to $\langle \Theta \rangle$, so that $\text{tr} \mathbf{C}_\Theta / n$ is the vanishing power of input fluctuations. Nevertheless, the approximation is the

same in both cases and reads

$$I(\Theta; \mathbf{R}) \approx \frac{\varepsilon^2}{2} \text{tr} [\mathbf{J}(\theta_0|\mathbf{R}) \mathbf{C}_{\tilde{\Theta}}] = \frac{1}{2} \text{tr} [\mathbf{J}(\theta_0|\mathbf{R}) \mathbf{C}_\Theta], \quad (24)$$

where $\langle \Theta \rangle = \theta_0$.

Eqns. (10) and (24) are identical, although the assumptions on Θ are different. Consider for example the memoryless channel with power constraint $P \geq \varepsilon^2$ on the input and $\langle \Theta \rangle = 0$, so that Eq. (24) can be written as

$$I(\Theta; R) \approx \frac{\varepsilon^2}{2} J(0|R). \quad (25)$$

The capacity is achieved by any distribution of inputs with power $P_\Theta = \varepsilon^2 = P$, for example by the discrete distribution from Eq. (16) with $\Delta\theta = \sqrt{P}$, or by the Gaussian distribution $\mathcal{N}(0, P)$. Specifically, it is well known that the capacity of a power-constrained linear additive white Gaussian noise (AWGN) channel is [1]

$$\mathcal{C} = \frac{1}{2} \ln \left(1 + \frac{P}{N} \right), \quad (26)$$

where P is the power constraint on the input and N is the noise power, and that the capacity is achieved by a normal distribution $\mathcal{N}(0, P)$. The signal-to-noise ratio (SNR) is then defined as $\text{SNR} = P/N$. By expanding Eq. (26) to first order in P for $P \ll N$ we have $\mathcal{C} \approx P/N/2$, which corresponds exactly to Eq. (25), since for the Gaussian additive noise holds $J(0|R) = 1/N$. A detailed review of AWGN channel capacity and its different approximations for different SNR regimes (including the high-noise approximation above) can be found in [37]. The conclusion that in the vanishing input-power limit the capacity of AWGN channel can be achieved by both discrete and $\mathcal{N}(0, P)$ distributions is not so surprising in the light of some recent research on the AWGN channels [38]. It has been shown, that although the optimal input distribution is generally $\mathcal{N}(0, P)$, the capacity can be near-achieved by a discrete distribution, and specially, if $P \ll N$ the other possible capacity-bearing distribution is indeed binary discrete. The methods employed in [38] are, however, different from our approach. We further discuss the compatibility of Eq. (24) with the exact results obtained for non-white AGN channels in the low-input power regime in the Results section of this paper.

C. Simple lower bound on memoryless channel capacity

We have demonstrated in the previous sections, that if the input to the memoryless channel is weak (in amplitude or power), the optimal distribution is discrete and binary. Therefore the channel capacity cannot be more than 1 bit. Note, however, that the capacity can be larger than 1 bit for channels with memory under certain circumstances, as we demonstrate in the Results section.

It follows from the proof in Appendix A, that the Fisher information arises in Eq. (10) from Taylor-expanding the involved KL distances in the expression for MI. More precise approximation to channel capacity, \mathcal{C}_{bin} , can be thus obtained without Taylor expansions, just by substituting the discrete input distribution from Eq. (16) into Eq. (2),

$$\mathcal{C}_{\text{bin}} = \frac{1}{2} D_{\text{KL}} [f(r|\theta_0 - \Delta\theta) \| p(r)] + \frac{1}{2} D_{\text{KL}} [f(r|\theta_0 + \Delta\theta) \| p(r)], \quad (27)$$

where $p(r) = f(r|\theta_0 - \Delta\theta)/2 + f(r|\theta_0 + \Delta\theta)/2$. The parameter $\Delta\theta$ is half of the maximum input amplitude for amplitude-constrained channels, and $\Delta\theta = \sqrt{P}$ for power-constrained channels.

Eq. (27) is the lower bound on the true capacity, $\mathcal{C} \geq \mathcal{C}_{\text{bin}}$, which holds whether the amplitude (or power) is small or not. The extension of Eq. (27) to channels with memory is not straightforward, for example the calculation of \mathcal{C}_{bin} would require numerical evaluation of possibly high-dimensional integrals which may not be numerically stable [39]. Therefore for channels with memory we propose to employ Eq. (10) as the simplest method.

IV. RESULTS FOR SELECTED SYSTEMS

A. Memoryless channels

1. Amplitude constrained linear AWGN channel

The capacity and capacity-bearing input distributions of the linear AWGN channel,

$$R = \Theta + Z, \quad (28)$$

where r.v. Z is zero-mean Gaussian and the input is constrained in amplitude, were studied in detail in [18]. Contrary to the well known Eq. (26) for the input power constrained channel, no closed-form expression for capacity exists in the amplitude constrained version, moreover the optimal input distribution is known to be discrete with finite set of mass points.

We assume $\theta_0 = 0$, the maximal input amplitude is $2\Delta\theta$, thus the input is bound to lie in the interval $[-\Delta\theta, \Delta\theta]$. Furthermore we assume that the power of the noise is $N = 1$, so the noise is described by the standard normal r.v., $Z \sim \mathcal{N}(0, 1)$. Eq. (15) then becomes

$$\mathcal{C}_{\text{high}} = \frac{1}{2} (\Delta\theta)^2. \quad (29)$$

The binary approximation, \mathcal{C}_{bin} given by Eq. (27), has to be evaluated numerically. Additionally, we also investigate the *low noise* approximation to MI, \mathcal{C}_{low} , which is also based on FI [12, 21, 22],

$$\mathcal{C}_{\text{low}} = \ln \frac{\int_{\Theta} \sqrt{J(\theta|R)} d\theta}{\sqrt{2\pi e}}. \quad (30)$$

Eq. (30) is a lower bound on the true channel capacity, $\mathcal{C} \geq \mathcal{C}_{\text{low}}$, tight with the vanishing noise in the information transmission. In the case of amplitude-constrained AWGN channel we have

$$\mathcal{C}_{\text{low}} = \ln \frac{2\Delta\theta}{\sqrt{2\pi e}}. \quad (31)$$

Fig. 2a. shows the comparison of the exact channel capacity (data taken from [16]) with $\mathcal{C}_{\text{high}}$, \mathcal{C}_{bin} and \mathcal{C}_{low} , expressed as functions of the signal-to-noise ratio (in dB), which is defined as [16]

$$\text{SNR} = 10 \log_{10} [(\Delta\theta)^2]. \quad (32)$$

The capacities are evaluated in bits which means converting the natural logarithms in Eqns. (15), (27) and (30) to base 2, i.e., to divide the values by $\ln 2$. While \mathcal{C}_{low} and $\mathcal{C}_{\text{high}}$ provide good approximations only for rather high and small SNR values, the \mathcal{C}_{bin} approximation gives good results even for intermediate SNR values. A similar figure with additional approximations for the classical AWGN channel capacity can be found in [37].

2. Temporal neuronal coding

Recently, the information capacity of a memoryless neuronal model has been analyzed in detail [17]. It is assumed, that the neuronal response R is the interval between two consequent action potentials. In agreement with some experimental observations [40–43], the response for each input follows the gamma distribution,

$$f(r|\theta) = \frac{r^{\kappa-1} \exp(-r/\theta)}{\theta^{\kappa} \Gamma(\kappa)}, \quad (33)$$

where the parameter θ is assumed to be the input (stimulus intensity). Based on further experimental observations [44], the input is constrained in amplitude, $5/\kappa \leq \theta \leq 50/\kappa$. The exact capacity was calculated numerically by Ikeda and Manton [17] for $0.75 \leq \kappa \leq 4.5$.

While \mathcal{C}_{bin} has to be evaluated numerically, for the high and low noise approximations we have

$$\mathcal{C}_{\text{high}} = \frac{81}{242} \kappa, \quad \mathcal{C}_{\text{low}} = \ln \frac{\sqrt{\kappa} \ln 10}{\sqrt{2\pi e}}. \quad (34)$$

The results are shown in Fig. 2b. For the investigated values of κ , both $\mathcal{C}_{\text{high}}$ and \mathcal{C}_{bin} approximations give better results than \mathcal{C}_{low} , which suggests that this particular case of temporal coding falls within the “high noise” category. Neuronal responses often vary substantially across identical stimulus trials, thus the highly noisy information transmission is not unusual as reported from experimental measurements [45]. A simple model of a stochastic resonance in an electrosensory neuron, subject to sub-threshold (i.e., very weak) stimulation [25, 46] has been analyzed by employing $\mathcal{C}_{\text{high}}$ recently [24].

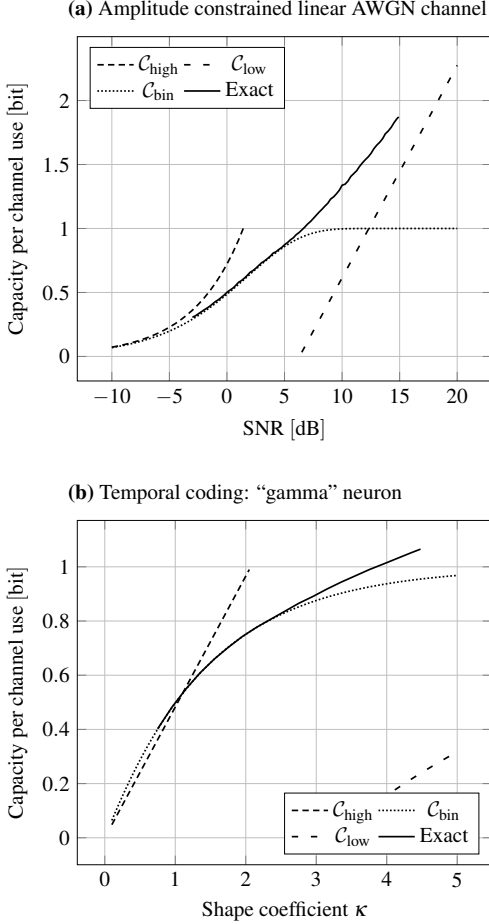


FIG. 2. Capacities and their approximations in memoryless channels. The high-noise capacity approximation (C_{high} , Eq. (15)) approximates the true capacity of the amplitude-constrained AWGN channel (a) well only for very low signal-to-noise ratios (SNR), just like the low-noise approximation (C_{low} , Eq. (30)) does for high SNRs. The binary-channel approximation (C_{bin} , Eq. (27)) holds well even for intermediate-low SNRs. The exact solution is taken from [16]. The information capacity of a simple model of neuronal coding (b) apparently falls into the high-noise category, since both C_{high} and C_{bin} approximate the true capacity (taken from [17]) better than C_{low} .

B. Linear Gaussian channel with memory and input power constraint

First, we demonstrate that Eq. (24) is compatible with exact results available on input power constrained linear AGN channels with memory [1, 19] in the limit of weak input power. The channel is defined as

$$\mathbf{R} = \boldsymbol{\Theta} + \mathbf{Z}, \quad (35)$$

where the zero-mean input is constrained in power P [1, p.277],

$$P \geq \frac{1}{n} \text{tr } \mathbf{C}_{\boldsymbol{\Theta}}, \quad (36)$$

and the noise is given by the multivariate normal distribution with covariance matrix $\mathbf{C}_{\mathbf{Z}}$, $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{Z}})$. The channel conditional p.d.f. is therefore

$$f(\mathbf{r}|\boldsymbol{\theta}) = \frac{1}{\sqrt{(2\pi)^n \det \mathbf{C}_{\mathbf{Z}}}} \exp [(\mathbf{r} - \boldsymbol{\theta})^T \mathbf{C}_{\mathbf{Z}}^{-1} (\mathbf{r} - \boldsymbol{\theta})], \quad (37)$$

and substituting Eq. (37) into Eq. (6) gives [27]

$$\mathbf{J}(\boldsymbol{\theta}|\mathbf{R}) = \mathbf{C}_{\mathbf{Z}}^{-1}, \quad (38)$$

which is independent of $\boldsymbol{\theta}$.

From the spectral decomposition theorem [34] follows that

$$\mathbf{C}_{\mathbf{Z}} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T, \quad (39)$$

where the matrix $\boldsymbol{\Lambda}$ is diagonal with positive elements and \mathbf{Q} is orthonormal. The capacity per channel use is then given by [19]

$$\mathcal{C} = \frac{1}{2n} \sum_{i=1}^n \ln \left(1 + \frac{m_i}{[\boldsymbol{\Lambda}]_{ii}} \right), \quad (40)$$

where the constants $m_i \geq 0$ are determined by the water-filling procedure [1, p.274], so that the power constraint given by Eq. (36) holds as $\sum_{i=1}^n m_i = nP$. Furthermore, the optimal input distribution is also multivariate normal, $\boldsymbol{\Theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\boldsymbol{\Theta}})$, with covariance matrix $\mathbf{C}_{\boldsymbol{\Theta}} = \mathbf{Q} \mathbf{M} \mathbf{Q}^T$ [19, p.279], where the diagonal matrix \mathbf{M} is defined as $[\mathbf{M}]_{ii} = m_i$.

In order to obtain the vanishing input power limit of Eq. (40), we observe that as $P \rightarrow 0$ also $m_i \rightarrow 0$, so we can expand Eq. (40) as

$$\mathcal{C} \approx \frac{1}{2n} \sum_{i=1}^n \frac{m_i}{[\boldsymbol{\Lambda}]_{ii}} = \frac{1}{2n} \text{tr} (\boldsymbol{\Lambda}^{-1} \mathbf{M}). \quad (41)$$

By combining Eqns. (38), (39), (41) and basic properties of matrix inverse and trace [34] we have

$$\begin{aligned} \mathcal{C} &\approx \frac{1}{2n} \text{tr} [(\mathbf{Q}^T \mathbf{C}_{\mathbf{Z}} \mathbf{Q})^{-1} \mathbf{M}] = \frac{1}{2n} \text{tr} [\mathbf{Q}^T \mathbf{C}_{\mathbf{Z}}^{-1} \mathbf{Q} \mathbf{M}] = \\ &= \frac{1}{2n} \text{tr} [\mathbf{C}_{\mathbf{Z}}^{-1} \mathbf{Q} \mathbf{M} \mathbf{Q}^T] = \frac{1}{2n} \text{tr} [\mathbf{J}(\boldsymbol{\theta}|\mathbf{R}) \mathbf{C}_{\boldsymbol{\Theta}}], \end{aligned} \quad (42)$$

which corresponds to the capacity per channel use as $n \rightarrow \infty$, due to Eq. (24), for power achieving input, $\text{tr } \mathbf{C}_{\boldsymbol{\Theta}}/n = P$.

Next, we illustrate Eq. (42) on two simple models of Gaussian noise with memory.

1. AR(1) noise

The channel is given by Eqns. (35) and (36), with Z_i 's following the AR(1) process: $Z_i = \varrho Z_{i-1} + X_i$, where $-1 < \varrho < 1$ is the correlation coefficient, $\varrho = \text{corr}(Z_i, Z_{i-1})$, and X_i are independently distributed

standard normal r.v.'s, $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ [32]. The noise covariance matrix has elements

$$[\mathbf{C}_Z]_{ik} = \varrho^{|i-k|}, \quad (43)$$

and its inverse, equal to the FI matrix by Eq. (38), is tridiagonal,

$$\mathbf{J}(\boldsymbol{\theta}|\mathbf{R}) = \frac{1}{1-\varrho^2} \begin{pmatrix} 1 & -\varrho & 0 & \cdots & 0 \\ -\varrho & 1+\varrho^2 & -\varrho & \cdots & 0 \\ 0 & -\varrho & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & 1+\varrho^2 & -\varrho \\ 0 & 0 & 0 & -\varrho & 1 \end{pmatrix}. \quad (44)$$

We denote the correlation coefficient between consequent inputs as $c = \text{corr}(\Theta_i, \Theta_{i+1})$. The MI per channel use for maximum power achieving input, $P = \text{tr} \mathbf{C}_\Theta / n$, can be found exactly by employing Eq. (24),

$$\lim_{n \rightarrow \infty} \frac{1}{n} I(\boldsymbol{\Theta}; \mathbf{R}) = \frac{P}{2} \frac{\varrho^2 + 1 - 2c\varrho}{1 - \varrho^2}. \quad (45)$$

For $\varrho = 0$ (memoryless channel) the value of c does not matter as discussed earlier. The capacity per channel use is

$$\mathcal{C}_{\text{high}} = \frac{P}{2} \frac{\varrho^2 + 1 + 2|\varrho|}{1 - \varrho^2}, \quad (46)$$

since $\sup_{-1 < c < 1} (-c\varrho) = |\varrho|$. The capacity in bits per vanishing input power, $\mathcal{C}_{\text{high}}/P$, is shown in Fig. 3 in dependence on the noise correlation ϱ . Note that from Eq. (46) follows $\mathcal{C}_{\text{high}}/P \rightarrow \infty$ as $|\varrho| \rightarrow 1$, i.e., as the noise correlation increases, its corrupting power decreases and in the limit we can approach the noiseless channel.

2. MA(1) noise

The channel is given by Eqns. (35) and (36), r.v.'s Z_i follow the MA(1) process, $Z_i = X_i - \gamma X_{i-1}$, where $-1 < \gamma < 1$ is the parameter of the process and $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. The parameter of the MA(1) process and the correlation coefficient $\varrho = \text{corr}(Z_i, Z_{i-1})$ are related as $\varrho = -\gamma/(1 + \gamma^2)$, and therefore $-0.5 < \varrho < 0.5$ [32]. The covariance matrix of the MA(1) process is tridiagonal, and its inverse has all elements non-zero, although decreasing in absolute value with the distance from the main diagonal, see Fig. 4a, b.

Recently, a closed form expression for \mathbf{C}_Z^{-1} of the MA(1) process has been published [47]. The expression is rather complicated and we cannot evaluate the analogous limit to Eq. (45) in a closed form. Nevertheless, we approximate the capacity per channel use by considering n high enough, and the closed form expression for the elements of the FI matrix allows us to avoid numerical issues when inverting the covariance matrix. The capacity per vanishing input power, $\mathcal{C}_{\text{high}}/P$, is shown in

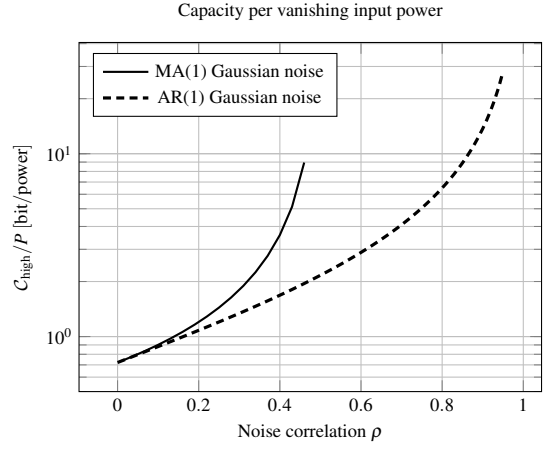


FIG. 3. The capacities per vanishing input powers for the AR(1) and MA(1) Gaussian additive noise models in dependence on the noise correlation coefficient ϱ (the graphs are symmetric in ϱ). Note that the capacity tends to infinity as $|\varrho| \rightarrow 0.5$ (the MA(1) model) and as $|\varrho| \rightarrow 1$ (the AR(1) model). In these limits, the corrupting power of the noise in the information transfer is decreased to the point, that the channel approaches the noiseless channel and the input value can be recovered perfectly.

Fig. 3. Note, that for $n \leq 2000$ we were unable to obtain stable values of $\mathcal{C}_{\text{high}}$ for $|\varrho| > 4.2$. This is caused by the fact, that the dominant terms of the FI matrix, and consequently $\mathcal{C}_{\text{high}}/P$, diverge to $+\infty$ as $|\varrho| \rightarrow 0.5$ (in a similar way as Eq. (46) does for $|\varrho| \rightarrow 1$). In other words, the dependence structure of the MA(1) process is sufficiently “rigid” even for intermediate correlation values, that by properly matching the input correlations we can approach the noiseless information transfer. The examples of optimally matched input signals are shown in Fig. 4c, d, e.

V. CONCLUSIONS

We derive approximate expression for mutual information in a broad class of discrete-time stationary channels (including those with memory) with continuous, but small, input. The input is restricted either in amplitude or in power and we study the optimality conditions on information transfer as the power or amplitude approach zero. We find that the input and channel properties are separated in the approximate formula, which allows us to study the optimality conditions in a convenient way. Specifically, we find that the increase of mutual information from zero power (or amplitude) for a given channel depends only on the input covariances.

For memoryless channels, the capacity cannot be more than 1 bit per channel use and the optimal input is unique discrete binary distribution in the small input amplitude case, but generally non-unique in the small input power

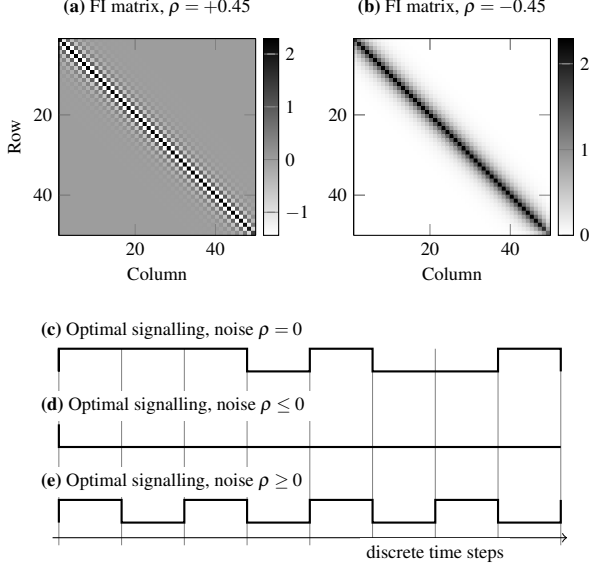


FIG. 4. Small input amplitude optimality conditions for linear channels with AR(1) or MA(1) additive Gaussian noise. The structure of the Fisher information matrix of the MA(1) model (panels (a) and (b) for $n = 50$) shows elements decaying in absolute value with distance from the main diagonal, sign changes occur for positively correlated MA(1) process and all elements are positive for $\rho \leq 0$. The structure of the FI matrix determines the covariance matrix of the optimal signal. Panel (c) shows the example optimal input to the memoryless channel (noise correlation $\rho = 0$): random switching between input values $+\sqrt{P}$ and $-\sqrt{P}$ (discrete binary input), where P is the input power constraint. The same capacity would be achieved by input values described by the normal distribution $\mathcal{N}(0, P)$, as discussed in the text. Depending on the sign of the noise correlation ρ , the optimal input is characterized by extremal value of correlation between consequent inputs (panels (d) and (e)). Note, that the capacity of the memoryless channel is achieved by (d) and (e) also, independently on the input correlations.

case. We demonstrate, that the effect of input correlations in memoryless channels is of higher order than the order of the capacity approximation, and thus the additional correlations do not decrease the capacity although they decrease the input entropy. We also provide a simple lower bound on capacity of memoryless channels subject to weak-stimulus constraints that gives better results in practical situations.

In channels with memory, the capacity can be greater than 1 bit and the input correlations play the most important role. We show, that the approximate formula includes the small input power limit of the exact solution for linear additive Gaussian noise channels with memory. We show, that by properly matching the input covariances to the dependence structure of the noise, we can approach in certain cases the noiseless channel even for intermediate values of the noise correlations.

ACKNOWLEDGMENTS

This work was supported by AV0Z50110509 and Centre for Neuroscience LC554. I thank Ales Nekvinda for helpful comments on the Appendix.

Appendix A: Capacity in the vanishing input amplitude

We introduce an auxiliary r.v. $\delta\Theta$ by employing Eq. (9) as

$$\delta\Theta = \Theta - \theta_0, \quad (\text{A1})$$

so that for all i holds $\delta\theta_i \in [-\Delta\theta, \Delta\theta]$. The p.d.f. of r.v. $\delta\Theta$ is denoted as $\pi(\delta\theta)$. Mutual information $I(\Theta; \mathbf{R})$ from Eq. (2) can be written in terms of r.v. $\delta\Theta$, whether $\|\Delta\theta\|$ is small or not as

$$I(\Theta; \mathbf{R}) = \langle D_{\text{KL}}[f(\mathbf{r}|\theta_0 + \delta\theta) \parallel \langle f(\mathbf{r}|\theta_0 + \delta\theta) \rangle_{\delta\theta}] \rangle_{\delta\theta}. \quad (\text{A2})$$

In order to approximate $I(\Theta; \mathbf{R})$ around θ_0 in terms of $\delta\theta$ for small $\|\Delta\theta\|$, we need to expand the KL distance in Eq. (A2). We introduce

$$\varphi(\mathbf{r}, \theta_0 + \delta\theta) = f(\mathbf{r}|\theta_0 + \delta\theta) \ln f(\mathbf{r}|\theta_0 + \delta\theta), \quad (\text{A3})$$

$$\psi(\mathbf{r}, \theta_0 + \delta\theta) = f(\mathbf{r}|\theta_0 + \delta\theta) \ln \langle f(\mathbf{r}|\theta_0 + \delta\theta) \rangle_{\delta\theta} \quad (\text{A4})$$

and rewrite the KL distance as

$$\begin{aligned} D_{\text{KL}}[f(\mathbf{r}|\theta_0 + \delta\theta) \parallel \langle f(\mathbf{r}|\theta_0 + \delta\theta) \rangle_{\delta\theta}] &= \\ &= \int_{\mathbf{R}} [\varphi(\mathbf{r}, \theta_0 + \delta\theta) - \psi(\mathbf{r}, \theta_0 + \delta\theta)] d\mathbf{r}, \end{aligned} \quad (\text{A5})$$

thus reducing the problem to expanding $\varphi(\mathbf{r}, \theta)$ and $\psi(\mathbf{r}, \theta)$. While the Taylor expansion of $\varphi(\mathbf{r}, \theta)$ is straightforward, the expansion of the logarithm of the expected value of $f(\mathbf{r}|\theta)$ in $\psi(\mathbf{r}, \theta)$ is examined in the following Lemma.

Lemma 1. *Let $f(\mathbf{r}|\theta)$ be twice continuously differentiable with respect to θ . Then for a chosen θ_0 , r.v. $\delta\Theta \sim \pi(\delta\theta)$ and $\Delta\theta$ such, that for all i holds $\Delta\theta > 0$ and $-\Delta\theta \leq \delta\theta_i \leq \Delta\theta$, there exists $P > 0$ such, that the following approximation for small enough $\|\Delta\theta\|$ holds*

$$\ln \langle f(\mathbf{r}|\theta_0 + \delta\theta) \rangle_{\delta\theta} \approx \ln f(\mathbf{r}|\theta_0) + \langle \delta\Theta \rangle_{\delta\theta}^T \frac{\nabla f(\mathbf{r}|\theta_0)}{f(\mathbf{r}|\theta_0)}, \quad (\text{A6})$$

where $\nabla f(\mathbf{r}|\theta_0) = \nabla f(\mathbf{r}|\theta)|_{\theta=\theta_0}$, the gradient is taken with respect to θ and $\langle \delta\Theta \rangle_{\delta\theta} = \langle \delta\theta \rangle_{\delta\theta}$ is the expectation of r.v. $\delta\Theta$. The maximum error of expansion (A6) is bounded by $P\|\Delta\theta\|^2$.

Proof. From the continuity of second derivatives of $f(\mathbf{r}|\theta)$ around θ_0 follows

$$\left| \frac{\partial^2 f(\mathbf{r}|\theta)}{\partial \theta_i \partial \theta_j} \right| \leq M, \quad (\text{A7})$$

for all i, j . The Taylor expansion of $f(\mathbf{r}|\boldsymbol{\theta})$ around $\boldsymbol{\theta}_0$ in terms of $\delta\boldsymbol{\theta}$ reads

$$f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) \approx f(\mathbf{r}|\boldsymbol{\theta}_0) + \delta\boldsymbol{\theta}^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0), \quad (\text{A8})$$

and furthermore

$$\begin{aligned} |f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) - f(\mathbf{r}|\boldsymbol{\theta}_0) - \delta\boldsymbol{\theta}^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)| &\leq \\ &\leq nM\|\delta\boldsymbol{\theta}\|^2 \leq C\|\Delta\boldsymbol{\theta}\|^2. \end{aligned} \quad (\text{A9})$$

By integrating the expansion (A8), i.e., by taking the expectation with respect to r.v. $\delta\boldsymbol{\Theta}$, and by employing inequality (A9) it can be established that

$$\begin{aligned} \left| \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}) - f(\mathbf{r}|\boldsymbol{\theta}_0) - \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0) \right| &= \\ = \left| \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) [f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) - f(\mathbf{r}|\boldsymbol{\theta}_0) - \delta\boldsymbol{\theta}^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)] d(\delta\boldsymbol{\theta}) \right| &\leq \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) C\|\Delta\boldsymbol{\theta}\|^2 d(\delta\boldsymbol{\theta}) = C\|\Delta\boldsymbol{\theta}\|^2, \end{aligned} \quad (\text{A10})$$

and therefore the following expansion holds

$$\int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}) \approx f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0), \quad (\text{A11})$$

with the maximum error of order $\|\Delta\boldsymbol{\theta}\|^2$. From the Lagrange mean value theorem follows, that for $A, B > 0$ holds

$$|\ln A - \ln B| \leq \frac{1}{\min(A, B)} |A - B|. \quad (\text{A12})$$

We set $A = \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta})$, $B = f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)$, and combine the inequalities (A10) and (A12) to obtain

$$\begin{aligned} |\ln A - \ln B| &= \left| \ln \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}) - \ln [f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)] \right| \leq \\ &\leq \frac{1}{\min(A, B)} |A - B| \leq \frac{1}{\min(A, B)} C\|\Delta\boldsymbol{\theta}\|^2, \end{aligned} \quad (\text{A13})$$

where $\min(A, B)$ is finite due to regularity of $f(\mathbf{r}|\boldsymbol{\theta})$. From the Taylor expansion of $\ln(a + x)$ around a in terms of x and the expression for the Lagrange remainder [48] we have

$$\left| \ln(a + x) - \ln(a) - \frac{x}{a} \right| \leq \frac{x^2}{a^2}. \quad (\text{A14})$$

Setting $a = f(\mathbf{r}|\boldsymbol{\theta}_0)$ and $x = \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)$ thus gives

$$\left| \ln [f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)] - \ln f(\mathbf{r}|\boldsymbol{\theta}_0) - \frac{\langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)}{f(\mathbf{r}|\boldsymbol{\theta}_0)} \right| \leq \frac{\|\nabla f(\mathbf{r}|\boldsymbol{\theta}_0)\|^2}{f^2(\mathbf{r}|\boldsymbol{\theta}_0)} \|\Delta\boldsymbol{\theta}\|^2. \quad (\text{A15})$$

Finally, we apply the triangle inequality for absolute value, $|\alpha - \beta| \leq |\alpha - \gamma| + |\gamma - \beta|$, setting

$$\alpha = \ln A = \ln \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}), \quad \beta = \ln f(\mathbf{r}|\boldsymbol{\theta}_0) + \frac{\langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)}{f(\mathbf{r}|\boldsymbol{\theta}_0)}, \quad (\text{A16})$$

$$\gamma = \ln B = \ln [f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)], \quad (\text{A17})$$

and by combining inequalities (A13) and (A15) we obtain

$$\begin{aligned}
& \left| \ln \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}) - \ln f(\mathbf{r}|\boldsymbol{\theta}_0) - \frac{\langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)}{f(\mathbf{r}|\boldsymbol{\theta}_0)} \right| \leq \\
& \leq \left| \ln \int_{\mathbf{R}} \pi(\delta\boldsymbol{\theta}) f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) d(\delta\boldsymbol{\theta}) - \ln \left[f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0) \right] \right| + \\
& + \left| \ln \left[f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0) \right] - \ln f(\mathbf{r}|\boldsymbol{\theta}_0) - \frac{\langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f(\mathbf{r}|\boldsymbol{\theta}_0)}{f(\mathbf{r}|\boldsymbol{\theta}_0)} \right| \leq \\
& \leq \frac{1}{\min(A, B)} C \|\Delta\boldsymbol{\theta}\|^2 + \frac{\|\nabla f(\mathbf{r}|\boldsymbol{\theta}_0)\|^2}{f^2(\mathbf{r}|\boldsymbol{\theta}_0)} \|\Delta\boldsymbol{\theta}\|^2 = P \|\Delta\boldsymbol{\theta}\|^2, \quad (\text{A18})
\end{aligned}$$

and therefore

$$\ln \langle f(\mathbf{r}|\boldsymbol{\theta}_0 + \delta\boldsymbol{\theta}) \rangle_{\delta\boldsymbol{\theta}} \approx \ln f(\mathbf{r}|\boldsymbol{\theta}_0) + \langle \delta\boldsymbol{\Theta} \rangle^\top \frac{\nabla f(\mathbf{r}|\boldsymbol{\theta}_0)}{f(\mathbf{r}|\boldsymbol{\theta}_0)}, \quad (\text{A19})$$

with error of order $\|\Delta\boldsymbol{\theta}\|^2$. \square

In the following we set $\varphi \equiv \varphi(\mathbf{r}, \boldsymbol{\theta}_0 + \delta\boldsymbol{\theta})$, $\psi \equiv \psi(\mathbf{r}, \boldsymbol{\theta}_0 + \delta\boldsymbol{\theta})$, $f \equiv f(\mathbf{r}|\boldsymbol{\theta}_0)$ and $\nabla f \equiv \nabla f(\mathbf{r}|\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$ for shorthand, and by repeatedly applying Lemma 1 and keeping in mind the rules for derivatives $(fg)'' = f''g + 2f'g' + fg''$, and $(\ln f)'' = f''/f - (f'/f)^2$, we obtain the expansions

$$\begin{aligned}
\varphi & \approx f \ln f + \delta\boldsymbol{\theta}^\top \ln f \nabla f + \delta\boldsymbol{\theta}^\top \nabla f + \\
& + \frac{1}{2} \delta\boldsymbol{\theta}^\top \ln f \nabla \nabla^\top f \delta\boldsymbol{\theta} + \delta\boldsymbol{\theta}^\top \frac{\nabla f \nabla^\top f}{f} \delta\boldsymbol{\theta} + \\
& + \frac{1}{2} \delta\boldsymbol{\theta}^\top f \left[\frac{\nabla \nabla^\top f}{f} - \frac{\nabla f \nabla^\top f}{f^2} \right] \delta\boldsymbol{\theta}, \quad (\text{A20})
\end{aligned}$$

$$\begin{aligned}
\psi & \approx f \ln f + \delta\boldsymbol{\theta}^\top \ln f \nabla f + \langle \delta\boldsymbol{\Theta} \rangle^\top \nabla f + \\
& + \frac{1}{2} \delta\boldsymbol{\theta}^\top \ln f \nabla \nabla^\top f \delta\boldsymbol{\theta} + \delta\boldsymbol{\theta}^\top \frac{\nabla f \nabla^\top f}{f} \langle \delta\boldsymbol{\Theta} \rangle + \\
& + \frac{1}{2} \langle \delta\boldsymbol{\Theta} \rangle^\top f \left[\frac{\nabla \nabla^\top f}{f} - \frac{\nabla f \nabla^\top f}{f^2} \right] \langle \delta\boldsymbol{\Theta} \rangle. \quad (\text{A21})
\end{aligned}$$

We substitute these expansions into Eq. (A5), and by

applying the regularity conditions (8) we have

$$\begin{aligned}
\int_{\mathbf{R}} [\varphi - \psi] d\mathbf{r} & \approx \frac{1}{2} \delta\boldsymbol{\theta}^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \delta\boldsymbol{\theta} - \\
& - \delta\boldsymbol{\theta}^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \langle \delta\boldsymbol{\Theta} \rangle + \frac{1}{2} \langle \delta\boldsymbol{\Theta} \rangle^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \langle \delta\boldsymbol{\Theta} \rangle, \quad (\text{A22})
\end{aligned}$$

where we employed the definition (6) of Fisher information matrix for $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) = \mathbf{J}(\boldsymbol{\theta}|\mathbf{R})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$. Due to symmetry $\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) = [\mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R})]^\top$ holds

$$\delta\boldsymbol{\theta}^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \langle \delta\boldsymbol{\Theta} \rangle = \frac{1}{2} \left[\delta\boldsymbol{\theta}^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \langle \delta\boldsymbol{\Theta} \rangle + \langle \delta\boldsymbol{\Theta} \rangle^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) \delta\boldsymbol{\theta} \right], \quad (\text{A23})$$

and so from Eq. (A2) we have

$$I(\boldsymbol{\Theta}; \mathbf{R}) \approx \frac{1}{2} \left\langle [\delta\boldsymbol{\theta} - \langle \delta\boldsymbol{\Theta} \rangle]^\top \mathbf{J}(\boldsymbol{\theta}_0|\mathbf{R}) [\delta\boldsymbol{\theta} - \langle \delta\boldsymbol{\Theta} \rangle] \right\rangle_{\delta\boldsymbol{\theta}}. \quad (\text{A24})$$

The covariance matrix $\mathbf{C}_{\delta\boldsymbol{\Theta}}$ of r.v. $\delta\boldsymbol{\Theta}$ is defined as

$$\mathbf{C}_{\delta\boldsymbol{\Theta}} = \left\langle [\delta\boldsymbol{\theta} - \langle \delta\boldsymbol{\Theta} \rangle] [\delta\boldsymbol{\theta} - \langle \delta\boldsymbol{\Theta} \rangle]^\top \right\rangle_{\delta\boldsymbol{\theta}}, \quad (\text{A25})$$

and obviously $\mathbf{C}_{\delta\boldsymbol{\Theta}} = \mathbf{C}_{\delta\boldsymbol{\Theta}}^\top$. Since $\boldsymbol{\theta}_0$ is fixed, and $\boldsymbol{\Theta} = \delta\boldsymbol{\Theta} + \boldsymbol{\theta}_0$, the covariance matrices of r.v. $\boldsymbol{\Theta}$ and r.v. $\delta\boldsymbol{\Theta}$ are equal, $\mathbf{C}_{\boldsymbol{\Theta}} = \mathbf{C}_{\delta\boldsymbol{\Theta}}$. Furthermore, the law of matrix multiplication gives $[\mathbf{A}\mathbf{B}]_{ik} = \sum_j [\mathbf{A}]_{ij} [\mathbf{B}]_{jk}$, thus summing along $i = k$ gives the trace, i.e., $\text{tr}(\mathbf{A}\mathbf{B}) = \sum_i [\mathbf{A}\mathbf{B}]_{ii} = \sum_{i,j} [\mathbf{A}]_{ij} [\mathbf{B}]_{ji}$. Therefore, Eq. (A24) can be written in a compact form as Eq. (10).

-
- [1] T. M. Cover and J. A. Thomas, *Elements of Information Theory* (John Wiley and Sons, Inc., New York, 1991).
 - [2] S. Verdú, IEEE Trans. Inf. Theory, **44**, 2057 (1998).
 - [3] C. H. Bennett and P. W. Shor, IEEE Trans. Inf. Theory, **44**, 2724 (1998).
 - [4] N. J. Cerf, J. Clavareau, C. Macchiavello, and J. Roland, Phys. Rev. A, **72**, 042330 (2005).
 - [5] R. G. Gallager, *Information theory and reliable communication* (John Wiley and Sons, Inc., New York, USA, 1968).

- [6] M. D. McDonnell and A. P. Flitney, Phys. Rev. E, **80**, 60102(R) (2009).
- [7] P. Mitra and J. B. Stark, Nature, **411**, 1027 (2001).
- [8] K. S. Turitsyn, S. A. Derevyanko, I. V. Yurkevich, and S. K. Turitsyn, Phys. Rev. Lett., **91**, 203901 (2003).
- [9] J. J. Atick, Network: Comput. Neural Syst., **3**, 213 (1992).
- [10] A. Borst and F. E. Theunissen, Nature Neurosci., **2**, 947 (1999).
- [11] S. B. Laughlin, Z. Naturforsch., **36**, 910 (1981).

- [12] M. D. McDonnell and N. G. Stocks, Phys. Rev. Lett., **101**, 058103 (2008).
- [13] L. Kostal, P. Lansky, and J.-P. Rospars, PLoS Comp. Biol., **4**, e1000053 (2008).
- [14] F. Rieke, R. de Ruyter van Steveninck, D. Warland, and W. Bialek, *Spikes: Exploring the Neural Code* (MIT Press, Cambridge, 1997).
- [15] R. Stein, Biophys. J., **7**, 797 (1967).
- [16] J. Dauwels, in *Proc. of the 26th Symposium on Inf. Theory in the BENELUX* (2005) pp. 221–228.
- [17] S. Ikeda and J. H. Manton, Neural Comput., **21**, 1714 (2009).
- [18] J. G. Smith, Inform. Control, **18**, 203 (1971).
- [19] R. W. Yeung, *Information theory and network coding* (Springer Verlag, Berlin, 2008).
- [20] R. J. McEliece, *The theory of information and coding* (Cambridge University Press, Cambridge, UK, 2002).
- [21] J. M. Bernardo, J. Roy. Stat. Soc. B, **41**, 113 (1979).
- [22] N. Brunel and J.-P. Nadal, Neural Comput., **10**, 1731 (1998).
- [23] B. S. Clarke and A. R. Barron, IEEE Trans. Inf. Theory, **36**, 453 (1990).
- [24] L. Kostal and P. Lansky, Phys. Rev. E, **81**, 050901(R) (2010).
- [25] P. E. Greenwood, L. M. Ward, D. F. Russell, A. Neiman, and F. Moss, Phys. Rev. Lett., **84**, 4773 (2000).
- [26] R. B. Ash, *Information Theory* (Dover, New York, 1965).
- [27] S. M. Kay, *Fundamentals of statistical signal processing: estimation theory* (Prentice Hall, New Jersey, 1993).
- [28] S. Kullback, *Information theory and statistics* (Dover, New York, 1968).
- [29] M. Salicrú and I. J. Taneja, Inform. Sciences, **72**, 251 (1993).
- [30] S. I. Amari and H. Nagaoka, *Methods of information geometry* (Amer. Math. Soc., Providence, USA, 2007).
- [31] J. J. Rissanen, IEEE Trans. Inf. Theory, **42**, 40 (1996).
- [32] M. Kendall, A. Stuart, and J. K. Ord, *The advanced theory of statistics. Vol. 3* (Charles Griffin, London, 1983).
- [33] Consider for example the matrix $\begin{pmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{pmatrix}$, which is positive-semidefinite, while $\begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix}$ is not. If the desired covariance matrix cannot be formed from Eq. (12), then C_{high} is less than as given by Eq. (13).
- [34] R. A. Horn and C. R. Johnson, *Matrix analysis* (Cambridge University Press, New York, 1985).
- [35] S. Verdú, IEEE Trans. Inf. Theory, **36**, 1019 (1990).
- [36] R. C. Dorf, *The Electrical Engineering Handbook* (CRC Press, Boca Raton, USA, 1997).
- [37] G. D. Forney and G. Ungerboeck, IEEE Trans. Inf. Theory, **44**, 2384 (1998).
- [38] J. Huang and S. P. Meyn, IEEE Trans. Inf. Theory, **51**, 2336 (2005).
- [39] D. R. Cox and N. Reid, Biometrika, **91**, 729 (2004).
- [40] M. W. Levine, Biol. Cybern., **65**, 459 (1991).
- [41] D. E. McKeegan, Brain Res., **929**, 48 (2002).
- [42] C. Pouzat and A. Chaffiol, J. Neurosci. Methods, **181**, 119 (2009).
- [43] G. N. Reeke and A. D. Coop, Neural Comput., **16**, 941 (2004).
- [44] S. Shinomoto, K. Shima, and J. Tanji, Neural Comput., **15**, 2823 (2003).
- [45] M. Carandini, PLoS Biology, **2**, 1483 (2004).
- [46] P. E. Greenwood and P. Lansky, Biol. Cybern., **92**, 199 (2005).
- [47] B. C. Sutradhar and P. Kumar, Appl. Math. Lett., **16**, 317 (2003).
- [48] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions, With Formulas, Graphs, and Mathematical Tables* (Dover, New York, 1965).